

1 In this case, users wait an average of only 122.5 milliseconds instead of 242.5, obtained  
2 using first-come, first served scheduling. More benefit is delivered sooner, and only less  
3 valuable and more costly jobs are delayed. Additionally, a pre-emptive scheduler has the  
4 ability, continuously, to insert higher priority jobs in front of lower priority ones. So, with the  
5 above scheduling and pre-emptive scheduling, the final, 200-millisecond job may run later  
6 than fourth, but it will not run sooner.

7  
8 The invention can employs well-known techniques for optimization and job scheduling.  
9 With the present invention known optimization and job scheduling techniques can be used  
10 to provide efficient Web and Internet servers, independent of the particular optimization or  
11 scheduling technique used  
12

13 Classification can be done as a mathematical function of known and estimated parameters.  
14 The above-mentioned benefit value may be a function of many parameters such as:  
15 requested URL, client IP address, cookie, login, connection quality and other distinguishing  
16 attributes. The cost may be more than just the time required to send or process the  
17 request, including, such factors as the required bandwidth, CPU, latency, data generation  
18 requirements, and total server load.

19  
20 In general, the benefit and cost are functions of known and estimated parameters can be  
21 described as follows:  
22

23 Benefit =  $B(p_1, p_2, p_3 \dots)$

24 Cost =  $B(p_1, p_2, p_3 \dots)$ , where  $p_N$  are known and estimated parameters.  
25

26 The functions may be tabular, or an actual mathematical function of the parameters or a  
27 combination of tabular and arithmetic functions. For example, a system can award points  
28 for a desirable URL request, one that is know to encourage commerce, compared to a  
29 reference or general information section of the site, which is read by both customers and  
30 non-customers.

31  
32 Order placement = 100 points

1 Browse catalog = 20 points

2 Download reference material = 10 points

3  
4 Points also go to requests that correspond to requestors who are, for example, good  
5 customers, as determined by the cookie, login, or, possibly, the IP source address.

6  
7 Returning customer = +20 points

8  
9 The connection quality determines how critical a connection is and how noticeable delays  
10 will be. This allows re-sequencing non-critical requests behind critical ones. In some  
11 cases, a modem user may not notice a slight delay but the DSL user will.

12  
13 RTT estimate < 100 ms, critical, speed sensitive

14 RTT estimate > 100 ms, non-critical

15  
16 The cost function is usually an estimate of the CPU required to generate the reply, the total  
17 time including latency required to generate the reply and the bandwidth required to send  
18 the reply back to the client. Other considerations include things such as the need for slots  
19 on application and database servers.

20  
21 For example, a typical response may cost 2 ms CPU, 20 ms latency, and 25k Bytes of data  
22 transfer. In general the optimal scheduler is one that delivers the maximum benefit, subject  
23 to the constraints that the total costs are less than the system limits in all cases.

24  
25 Example Two: Server connection and load management: A special case of a managed  
26 resource is a secondary server, usually an application server or database server. The  
27 server may suffer performance problems if its load is too high or if there are too many  
28 connections to the server from remote clients.

29  
30 Typical server response time vs. load

31 1 pending requests -- 10 ms avg. response time

32 5 pending requests -- 12 ms avg. response time

- 1 10 pending requests – 20 ms avg. response time
- 2 20 pending requests -- 50 ms avg. response time
- 3 100 pending requests – 400 ms avg. response time

4

5 Here we see that, initially, efficiency increases due to increased concurrency and

6 overlapping of requests that have both a latency (delay) and a processing (CPU)

7 component. After a certain point, the server becomes less efficient due to overhead of

8 maintaining many pending requests. Many application and database servers use operating

9 system threads or processes to handle simultaneous tasks. This results in diminishing

10 returns as threads corresponding to pending requests compete for synchronization

11 primitives and as the operating system is forced to switch back and forth among the

12 outstanding tasks.

13

14 In the above example, the server runs most efficiently at a load of around 10 pending

15 requests, 20 ms average response time, for a total of around 500 "hits" per second. If the

16 load is 100, with an average response time of 400 ms, then the throughput is about 250 hits

17 per second. Intelligent load management would maintain the load on the server at 10,

18 while queuing the remaining requests. As described in Example 1, this queuing has the

19 added benefit of being able to order or prioritize the requests within the queue, with

20 additional gains in throughput and reduction of average response time.

21

22 With the present invention requests can be handled by an intermediate server/optimizer,

23 which queues the connection and transfers the data back and forth between the requesting

24 client and the origin and application servers.

25

26 Net result due to intelligent load management, with 100 pending requests:

27 500 hits/sec, with an average response time of  $20 + (2\text{ms}) \cdot 90 = 200 \text{ ms}$ , compared with 250

28 hits/sec and 400 ms average response time with the standard application server.

29

30 Often, simply connecting to application and database servers slows the progress of tasks

31 executing on the server. In this case, it is helpful to off-load idle connections to an

32 intermediate server, which handles connection with an efficient queuing and I/O system.

A program flow diagram illustrating the operation of the system is given in Figure 2. First as indicated by block 201, the owner establishes a set of classifications priorities. These are stored in data base 123. If such a set are not as yet established the system utilizes a default set of priorities and classifications. The system receives requests from clients 100A to 100Z as indicated by block 203.

As indicated by block 295, the requests are classified in accordance with the classifications established by the system owner and stored in data base 123. Next as indicated by block 207, each request is prioritized in accordance with its classifications. Finally as indicated by block 208, the requests are scheduled in accordance with their priorities. Naturally higher priority requests are scheduled to be processes before lower priority requests.

Finally as indicated by block 209 the requests are sent to the web application server 110 and any outgoing traffic is sent to the network interface 120 for dispatch to the client machines 100A to 100Z.

Program 124 schedules the requests according to their priority and then prioritizes the requests in buffer 125. The requests are sent to the web application server 110 in accordance with their priority. A low priority requests which arrived first may reach web application server 112 after a high priority request which the system received at a later time.

In summary, the system includes programs that classify and priorities requests according to parameters established by the system operator. Different types of requests are provided with different priorities such that high priority requests are acted upon by the web application server 110 before low priority requests. This gives the system operator a great deal of flexibility in arranging the system provide a desired type of performance.

While in the embodiment described above, the classification, prioritization, and scheduling are done by three separate program routines, it should be understand, that the present invention relates to a program and system that performs this combination of functions. Those skilled in the art will realize that these functions can be performed using a wide variety of programming arrangements other than three separate programs.



unreachable". The kernel layer is designed to conserve system resources, especially RAM, so that the server will function optimally and degrade gracefully. Most systems exhibit non-linear delay vs. load characteristics, with a sharp knee in the curve at a critical load, indicating non-graceful degradation. The present invention will extend the curve by deferring lower priority and "housekeeping" tasks and by using system resources more efficiently.

The cost-benefit model enables prioritization of requests by content (URL) or requester (IP address, login, cookie), as well as according to more automatic criteria such as content length, resource requirements, or end-to-end connection quality. For example, if the server has one large request and ten small requests it may wish to service the ten small requests first, satisfying ten users, while adding an acceptable delay to the large request. Furthermore, shaving 100 milliseconds off a 200 millisecond RTT (round trip time) task would result in noticeable increase in interactivity. However, shaving 100 milliseconds off a 600-millisecond modem connection would not even be appreciated. This targeted, fine-grain optimization is enabled by characterizing the requests, estimating their resource requirements, then queuing up the required tasks in the correct sequence to optimize the model.

A server in accordance with the present invention improves efficiency by performing resource allocation and scheduling according to a cost-benefit model that is established both automatically and by the server's administrator. Requests to the server may be classified according to URL, URL parameters, requestor, connection quality, content size and generation requirements, server required, time of day, or any other identifying characteristic.

Each class of request may have its own priority level, benefit, maximum or weighted proportional share of total bandwidth, maximum or weighted proportional share of an assigned CPU, or priority of access to any system resource, server, or storage device.

Each class may have deadlines or constraints on delivery, with variable penalties for lateness and dropping. External, user objectives and constraints are translatable to internal

constraints, which determine CPU and bandwidth scheduling and proportional share at the segment (packet) granularity.

With the present invention, a custom stack can schedules packet (segment) departure based on deadlines and lateness penalties that have been established by the scheduler and allocator

The invention utilizes an event-driven framework. A custom OS layer which manages the classification, prioritization and scheduling reduces the overhead of the general-purpose OS in the server, and it provides better communication between multiple, simultaneous tasks that are in progress. The custom OS layer maps multiple request/reply tasks to fewer threads or a single thread of the host OS.

The custom OS layer uses knowledge of continuations and non-blocking activities, co-operative multi-tasking based on a trust relationship between tasks is enabled. This is similar to the technique that is often used in the design of simulators. Such techniques reduce the overhead of multiprogramming a large number of independent tasks. Monitoring or a "pulse function" detects blocked or deadlocked processes to transfer workload to other, functioning processes, of which, new server OS instances may be added as needed.

A customized protocol stack can reduce the cost of open connections that have no assigned or discernible pending tasks. Such connections store only a source address and port without the usual socket resources. This stack may run in parallel with the existing TCP/IP stack by intercepting relevant ports, protocols, and URL requests at the kernel level, affording them special treatment.

The present invention eliminates network layer overhead. The custom stack also provides feedback to the scheduler and allocator regarding connection quality and available TCP and IP resources. Similarly, the custom stack greatly reduces the cost of servicing in-memory "cached" replies by forgoing the need for creating "socket" resources to grant access of kernel data to user spaces.

The custom stack may multiplex multiple connections and data transfers “on the wire” and at the network layer into a single user-level connection at the OS-user/application layer.

With the present invention, customized applications replace the layered network interface with interprocess communication (i.e. IPC) and remote procedure calls (i. e. RPC) to communicate directly with the system more efficiently.

One valuable by-product of peeking into the TCP layer to glean connection information is the ability to provide the server owner with more detailed traffic statistics. The server statistics and quality of achieving the user-defined and automatic criteria is fed-back to a monitoring and reporting application, which then displays said information to the administrator.

In conclusion: The present invention provides an event-driven custom kernel for a server. The custom kernel provides scheduling and resource allocation. The custom kernel operates in accordance with a cost-benefit model which is optimized by request and task prioritization. The cost-benefit model prioritizes requests by content (URL or requestor (IP address, login, or cookies) as well as according to criteria such as content length, resource requirements, or end-to-end connection quality. Tasks are classified and prioritized before being run on the CPU. Bandwidth is regulated and data departure is scheduled according to task and server specific criteria that can be established by a user. Fine-grain optimization is achieved by characterizing the requests, estimating their resource requirements, then queuing up the required tasks in the correct sequence to optimize the model. Several types of Interprocess communication (i.e. IPC) and remote procedure calls (i. e. RPC) are used to efficiently communicate directly with the system. These include providing feedback information between layers, and sending data directly from an internal layer to a receiver that is not an adjacent layer via inter-process communication. Since the kernel in the server obtains information from the TCP and IP layers, detailed traffic statistics can be provided to the server owner.

While the invention has been shown and described with respect to preferred embodiments thereof, it will be understood by those skilled in the art, the various changes in form and



1 detail can be made without departing from the spirit and scope of the present invention.

2 The invention is limited only by the appended claims.

3

4 I claim:

5

[illegible]